# 2022 AI HW Summit:
# The Annual Top 10 List

Karl Freund

Founder and Principal Analyst

Cambrian-AI Research LLC

Karl.Freund@Cambrian-ai.com

@karlfreund

# The 2021 Top 10 List
➔ All are still Evolving

1. MLPerf Benchmarks
2. Synopsys DSO.ai
3. Sambanova, Groq, & Tenstorrent
4. AWS: Inferentia
5. Google: TPU-V4
6. Intel Habana Labs Gaudi
7. Graphcore 2$^{nd}$ Generation and the IPU-Machine
8. NVIDIA Grace
9. Cerebras WSE-2 and Brain-Scale AI
10. …. (The founding or Cambrian-AI, of course!)

# Why rob Banks? That's where the money is.
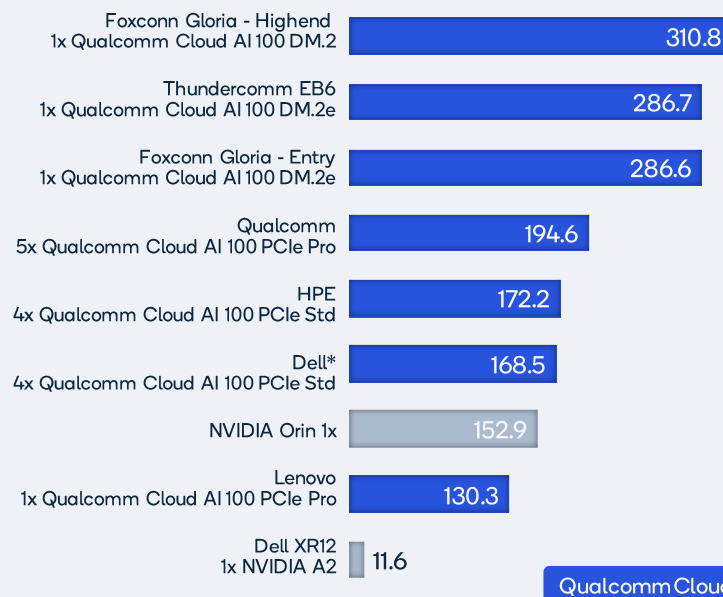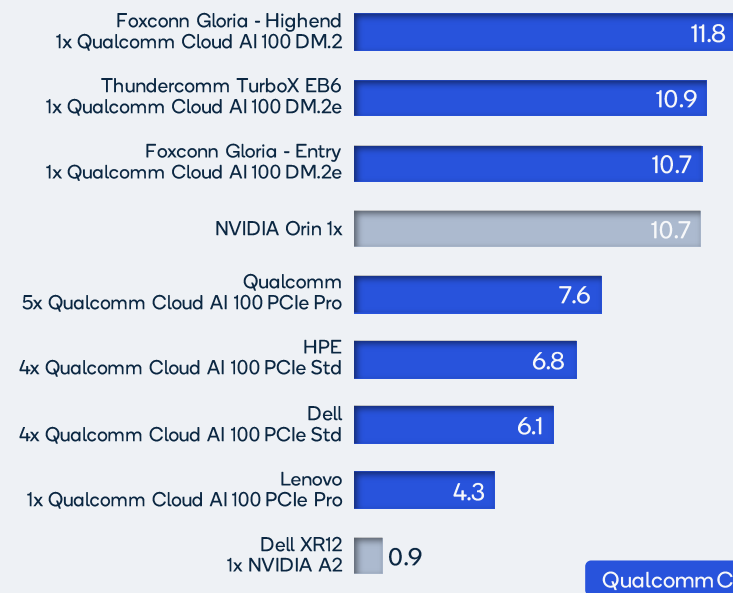
(C) Cambrian AI Research LLC

# MLPerf Inference 2.1

# Qualcomm Leads MLPerf Efficiency
## Over 200 submissions with Server Ecosystem

### Closed Edge Power - ResNet50
#### (Queries Per Sec/Watt)

| System | QPS/Watt |
|---|---|
| Foxconn Gloria - Highend<br>1x Qualcomm Cloud AI 100 DM.2 | 310.8 |
| Thundercomm EB6<br>1x Qualcomm Cloud AI 100 DM.2e | 286.7 |
| Foxconn Gloria - Entry<br>1x Qualcomm Cloud AI 100 DM.2e | 286.6 |
| Qualcomm<br>5x Qualcomm Cloud AI 100 PCIe Pro | 194.6 |
| HPE<br>4x Qualcomm Cloud AI 100 PCIe Std | 172.2 |
| Dell*<br>4x Qualcomm Cloud AI 100 PCIe Std | 168.5 |
| NVIDIA Orin 1x | 152.9 |
| Lenovo<br>1x Qualcomm Cloud AI 100 PCIe Pro | 130.3 |
| Dell XR12<br>1x NVIDIA A2 | 11.6 |

Qualcomm Cloud AI 100

### Closed Edge Power - BERT-99
#### (Queries Per Sec/Watt)

| System | QPS/Watt |
|---|---|
| Foxconn Gloria - Highend<br>1x Qualcomm Cloud AI 100 DM.2 | 11.8 |
| Thundercomm TurboX EB6<br>1x Qualcomm Cloud AI 100 DM.2e | 10.9 |
| Foxconn Gloria - Entry<br>1x Qualcomm Cloud AI 100 DM.2e | 10.7 |
| NVIDIA Orin 1x | 10.7 |
| Qualcomm<br>5x Qualcomm Cloud AI 100 PCIe Pro | 7.6 |
| HPE<br>4x Qualcomm Cloud AI 100 PCIe Std | 6.8 |
| Dell<br>4x Qualcomm Cloud AI 100 PCIe Std | 6.1 |
| Lenovo<br>1x Qualcomm Cloud AI 100 PCIe Pro | 4.3 |
| Dell XR12<br>1x NVIDIA A2 | 0.9 |

Qualcomm Cloud AI 100

## Most Power Efficient AI Edge Solution – Power Efficiency

MLPerf™ v2.1 submission IDs: Qualcomm 2.1-0105, 2.1-0104, 2.1-0103, 2.1-0108, HPE 2.1-0054 Lenovo 2.1-0081, Dell 2.1-0011, 2.1-0017, Nvidia 2.1-0096
Power efficiency in not MLPerf™ Metric. It is derived from MLPerf™ 2.1 Closed Edge Power submission.
* Dell is Closed Preview submission

6

# The Top 10 AI HW Innovations of 2022

1. **GrAI Matter Labs** – High Fidelity Edge processing
2. **D-Matrix:** In-Memory Computation
3. **Untether.AI:** At-Memory Computation
4. **Mythic:** Here comes Analog compute
5. **Esperanto** : 1000 RISC-V cores
6. **SimaAI:** MLSoC provides SW-Centric ML
7. **AMD**: Massive FLOPS w/ A100-class TOPS
8. **Graphcore**: Wafer on Wafer, and Good Computer
9. **Intel Habana Labs**: Gaudi 2 doubles NVIDIA A100 throughput
10. **NVIDIA** Hopper and Grace **Superchips**: the Next Generation of Compute
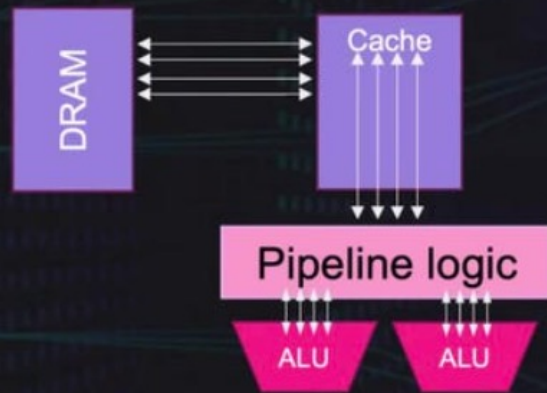
# GrAI Matter Labs:
## High Fidelity real-time inference

# D-Matrix: In-Memory Compute

# Untether AI: 30 TFLOPS/Watt



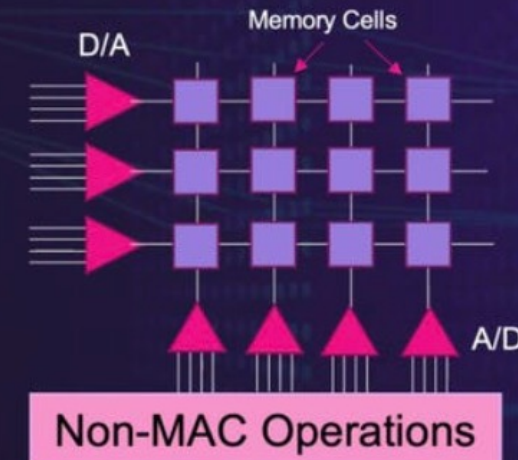At-Memory Compute Is the Sweet Spot for AI Acceleration

# Mythic: Here comes Analog Compute

# Esperanto: 1000 RISC-V Cores



## ML Recommendation performance per card comparisons
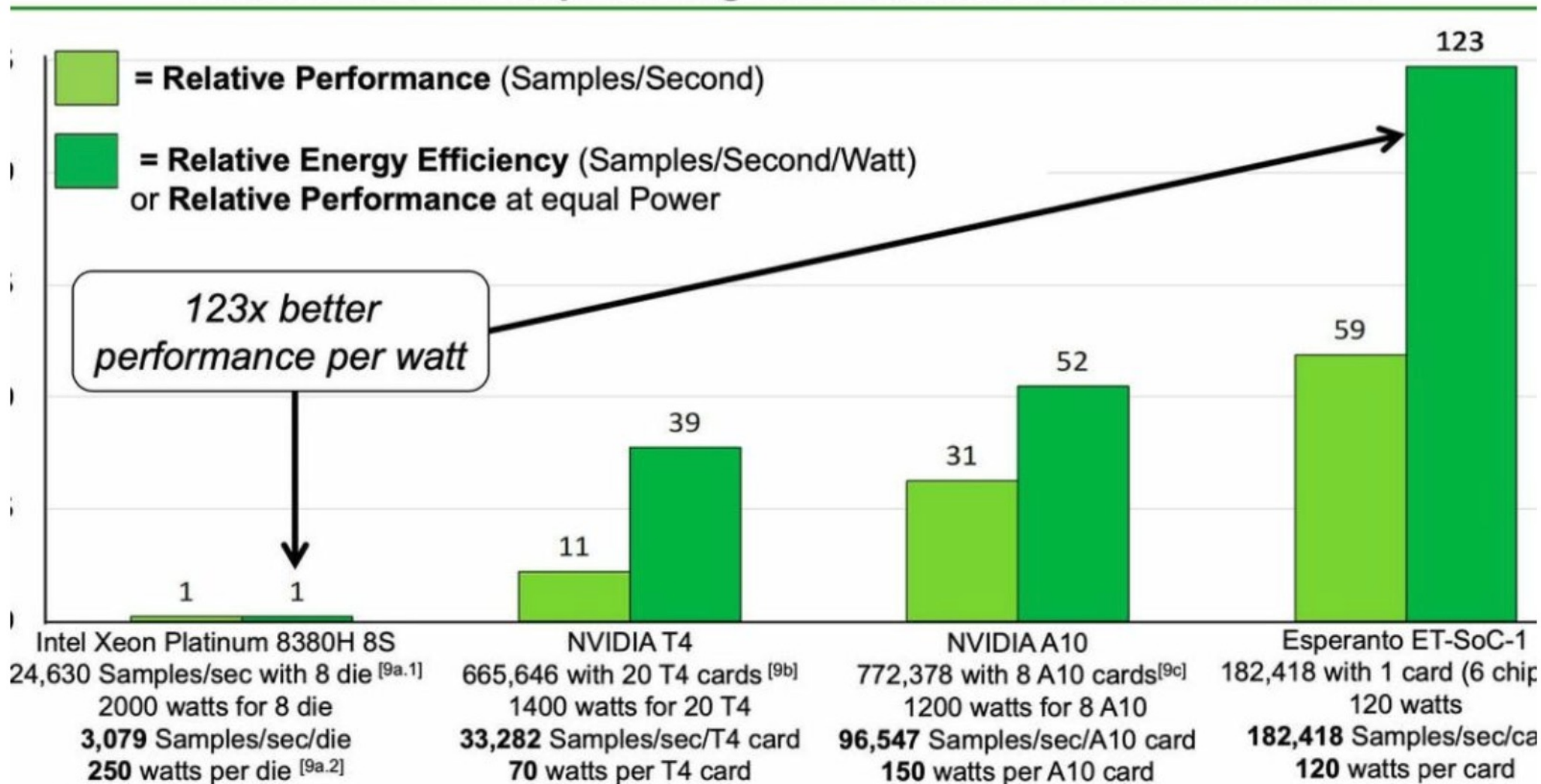### Based on MLPerf Deep Learning Recommendation Model benchmark [8]

■ = Relative Performance (Samples/Second)

■ = Relative Energy Efficiency (Samples/Second/Watt)
or Relative Performance at equal Power

123x better performance per watt

123

59

52

39

31

11

1    1

**Intel Xeon Platinum 8380H 8S**
24,630 Samples/sec with 8 die [9a.1]
2000 watts for 8 die
3,079 Samples/sec/die
250 watts per die [9a.2]

**NVIDIA T4**
665,646 with 20 T4 cards [9b]
1400 watts for 20 T4
33,282 Samples/sec/T4 card
70 watts per T4 card

**NVIDIA A10**
772,378 with 8 A10 cards [9c]
1200 watts for 8 A10
96,547 Samples/sec/A10 card
150 watts per A10 card

**Esperanto ET-SoC-1**
182,418 with 1 card (6 chip
120 watts
182,418 Samples/sec/ca
120 watts per card

# SiMa.ai MLSoC

# AMD Instinct MI200



## SHATTERING PERFORMANCE BARRIERS IN HPC & AI

| PEAK PERFORMANCE | A100 | MI200* | INSTINCT™ ADVANTAGE |
|---|---|---|---|
| FP64 VECTOR | 9.7 TF | 47.9 TF | 4.9X |
| FP32 VECTOR | 19.5 TF | 47.9 TF | 2.5X |
| FP64 MATRIX | 19.5 TF | 95.7 TF | 4.9X |
| FP32 MATRIX | N/A | 95.7 TF | N/A |
| FP16, BF16 MATRIX | 312 TF | 383 TF | 1.2X |
| MEMORY SIZE | 80 GB | 128 GB | 1.6X |
| MEMORY BANDWIDTH | 2.0 TB/s | 3.2 TB/s | 1.6X |

NOTE: THE A100 TF32 DATA FORMAT IS NOT IEEE FP32 COMPLIANT , SO NOT INCLUDED IN THIS COMPARISON.

*MI250x, SEE ENDNOTES: MI200-01, MI200-07
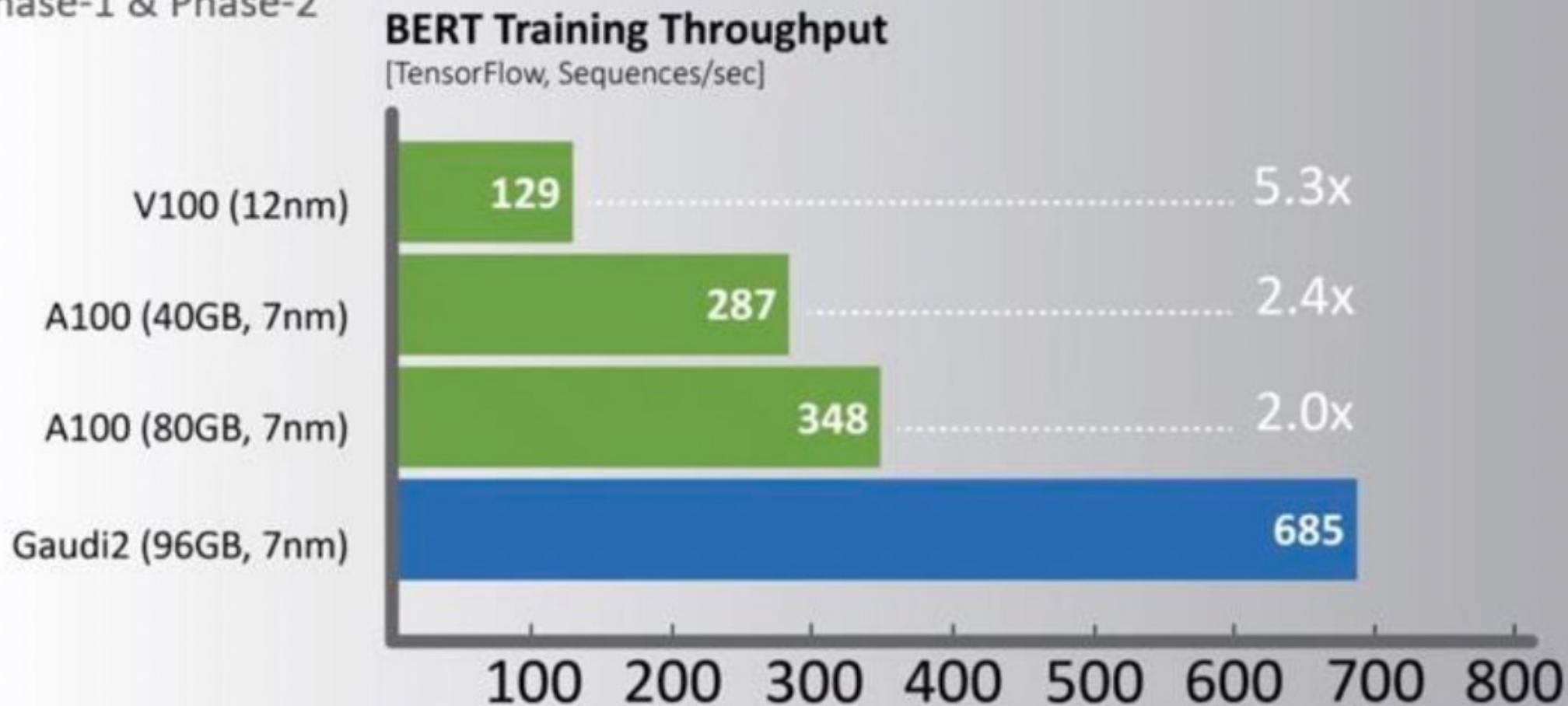
# Graphcore: Thinking BIG

# Intel Habana Labs Gaudi2: 2X A100



**Effective Throughput:**
Combining Phase-1 & Phase-2

**BERT Training Throughput**
[TensorFlow, Sequences/sec]

| | | |
|---|---|---|
| V100 (12nm) | 129 | 5.3x |
| A100 (40GB, 7nm) | 287 | 2.4x |
| A100 (80GB, 7nm) | 348 | 2.0x |
| Gaudi2 (96GB, 7nm) | 685 | |

100  200  300  400  500  600  700  800

# Grace Superchips:
# The Future of High Performance

# Closing Thoughts

- NVIDIA is practically unassailable in data center training, where Intel, Cerebras, Graphcore, SambaNova, Tenstorrent, and Groq are all attacking.
  - So, if you want to go after this space, you need <u>magic tech and a lot of money</u>

- The Edge market(s) are ripe with opportunities with <u>many niches</u>.

- Software matters even more than you think
  - NVIDIA improved Jetson by **50%** last MLPerf round

- BIG networks (LLMs) will eat the world
  - NVIDIA
  - Graphcore
  - Cerebras

# LLMs being used to create images from text

*Prompt: "oil on canvas painting + romanticism + landscape + a hay wain pulled by two horses as it crosses a river + a backdrop of mountains, trees, and clouds in the background + simple and idyllic depiction of rural life in England"*



Left: The Hay Wain by John Constable (Public Domain). Right: Image by Alberto Romero via Midjourney

# THANK YOU!

## Have a great conference!

Our clients include:



Come visit our site for news and research reports
http://www.Cambrian-ai.com